

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
2009-2010	L3	ELIZA en français	<p>On vous propose d'implémenter une version française de ELIZA. ELIZA est historiquement le premier succès spectaculaire du TAL. Il s'agit d'un petit agent conversationnel qui simule une conversation entre vous et un psychanalyste de patience infinie. Exemple de conversation en français :</p> <p>Système : "Comment allez-vous ?"</p> <p>Utilisateur : "Je suis déprimé."</p> <p>S : "Dites-moi pourquoi vous êtes déprimé."</p> <p>U : "Je passe trop de temps avec les ordinateurs."</p> <p>S : "Pourquoi mentionnez-vous les ordinateurs ?"</p> <p>U : "Parce que tu en es un."</p> <p>S : "Quels sont vos rapports avec les ordinateurs ?"</p> <p>U : "Difficiles."</p> <p>S : etc.</p> <p>Ce petit agent fonctionne principalement à l'aide de techniques élémentaires de pattern matching (expressions régulières) et vous demande de manipuler une pile simulant sa mémoire. Le sujet se focalisera sur les points suivants :</p> <ul style="list-style-type: none"> • Identifier un mot clé et axer la réponse sur ce mot clé (exemple : retenir un groupe de deux mots commençant par un déterminant, "les ordinateurs" dans "Je passe trop de temps avec les ordinateurs"), • Revenir à un mot clé quand l'énoncé courant ne permet de générer aucune réponse, • Exploiter des phrases à trou (exemple : "Pourquoi mentionnez-vous" + mot clé), • Traduire les mots référant aux interlocuteurs : "je", "vous", "moi", "me", etc. (exemple : "je suis déprimé" devient "vous êtes déprimé" dans "Dites-moi pourquoi vous êtes déprimé"). • Sortir des réponses types au hasard (exemple : "Comment allez-vous ?"). 	Java	Grégoire Winterstein	Difficile	2 personnes
		Détection de la langue d'un texte	<p>Pour détecter la langue d'un texte, on peut constituer une base de connaissances à partir d'un corpus de textes classés par langue. Pour chaque langue, un premier programme (à écrire) recueillera des statistiques significatives, basées sur les lettres (par exemple, il y a plus de "w" en anglais qu'en français). Le choix du modèle probabiliste employé et de ses paramètres (bigrammes, trigrammes) devra être justifié dans le rapport. On peut en proposer plusieurs et discuter de leurs avantages et inconvénients (rapport entre précision de la reconnaissance et volume de la base de connaissances ou longueur du texte nécessaire pour reconnaître sa langue).</p> <p>Un deuxième programme (à écrire), utilisera ces bases de connaissances pour reconnaître la langue d'un texte.</p> <p>Vous utiliserez un corpus d'apprentissage comprenant des textes plus ou moins variés d'un certain nombre de langues (le plus de langues possible, au moins 4). Le corpus de test ne devra pas contenir de texte appartenant au corpus d'apprentissage.</p> <p>Amélioration possible : gérer des textes dans différents encodages, pour les langues à alphabet non latin.</p>	Java	Grégoire Winterstein	Facile à moyen	2 personnes
		Correcteur orthographique	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> • de choisir un de ces remplaçants <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences • de fournir un remplaçant <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences o si le remplaçant n'est pas dans le lexique, 	Java	Grégoire Winterstein		2 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<input type="checkbox"/> de l'insérer dans le lexique pour ce texte, <input type="checkbox"/> de l'insérer dans le lexique stable <ul style="list-style-type: none"> • d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif. Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot : <ul style="list-style-type: none"> • repérage des bigrammes ou trigrammes impossibles en français • redoublement ou dé-doublement de consonnes • suppression/insertion de diacritiques • distance d'édition • etc Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.				
		Réaccentuation	Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre. On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple : La ou le français n'est pas accentué, il y a de la gene, mais quand le système m'accentue, je suis moins gene!	Java	Grégoire Winterstein		2 personnes
		Entités nommées	On regroupe sous le terme "entités nommées" les noms de personnes, de lieux, de dates, noms d'entreprises, adresses, etc. Il s'agit d'expressions qui dénotent une entité unique de façon presque indépendante du contexte. On s'intéresse aux entités nommées pour plusieurs raisons : <ul style="list-style-type: none"> • elles constituent des syntagmes qui peuvent être relativement complexes au point de vue syntaxique (par exemple une adresse, ou un nom d'association) dont le repérage préalable peut grandement simplifier une analyse syntaxique ; • dans une perspective de recherche d'information, la reconnaissance des entités nommées permet de savoir de quoi parle un texte ; • elles sont nécessaires pour la résolution des anaphores. Il s'agit dans ce projet de repérer de la façon la plus complète possible dans un texte étiqueté ou non, les entités de type "personne". Pour cela, on envisagera un algorithme en deux étapes (qui peuvent se répéter) : <ul style="list-style-type: none"> • au moyen de règles générales et de dictionnaires spécialisés (noms propres, amorces --- c'est-à-dire mots qui introduisent systématiquement des entités nommées, comme 'Melle', etc.), constitution d'une "table des symboles" des entités présentes dans le texte ; • à partir de cette table des symboles, et en tenant compte des formes variées sous lesquelles une même entité peut être désignée, recherche de nouvelles entités, voire de nouvelles règles trouvées précédemment. L'idée est que le programme s'enrichit au fur et à mesure qu'il est utilisé.	Java	Grégoire Winterstein	Facile à très difficile	2 ou 3 personnes
		Extraction de collocations en corpus	Les collocations sont des ensembles de deux ou plusieurs mots qui sont plus fréquemment cooccurrents que la normale. Il s'agit de constructions qui ont une syntaxe classique mais une sémantique qui n'est pas complètement compositionnelle.	Java, python ou perl et initiation à R (facultatif)	Grégoire Winterstein	Moyen	2 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>Par exemple on considère que thé fort est une collocation car ces deux mots apparaissent assez souvent en même temps et le sens de fort dans cette séquence est légèrement altéré : il ne s'agit pas d'une grande force physique mais plutôt d'une forte concentration d'un agent du thé.</p> <p>Ce sujet propose de partir à la découverte de différentes méthodes connues de détection de collocations en corpus reposant sur des tests d'hypothèses statistiques (comme par exemple le test du Chi2 ou le test exact de Fisher) en utilisant des logiciels d'analyse de données appropriés.</p>				
		Clavier téléphone (T9)	<p>La plupart des téléphones associent à chaque touche un certain nombre de lettres, permettant ainsi de transmettre des messages. Du fait qu'un chiffre ne correspond pas à une seule lettre, une suite de chiffres peut être ambiguë (elle peut correspondre à plus d'un mot). La suite 7-6-8-7, par exemple, correspond aux mots pour, sous et pots. Le but du projet est de réaliser un programme qui prend en entrée une suite de chiffre et un dictionnaire de formes fléchies et qui propose en sortie la liste des mots du dictionnaire qui correspondent à la suite des chiffres. Ces mots seront classés selon une fréquence représentée dans le dictionnaire. On attachera un soin particulier au choix de la structure de données, et à l'algorithme, qui doit faire le calcul en un temps raisonnable.</p>	Java, python ou perl	Grégoire Winterstein	Moyen	2 personnes
	M1	Marquage des antécédents pronominaux en corpus	<p>Il s'agit de mettre en oeuvre un ou plusieurs algorithmes de résolution d'anaphore sur un texte analysé syntaxiquement (Le French Treebank du laboratoire LLF pourra être utilisé).</p> <p>A partir des occurrences pronominales de 3e personne marquées pour le genre et le nombre de leurs antécédents, il s'agit de proposer un ou plusieurs candidats dans la même phrase (pour les relatifs et les réfléchis) ou non (pour les personnels).</p> <p>Les mots antécédents seront repérés par leur numéro dans la phrase (ID) éventuellement préfixé par le numéro de la phrase.</p>	Java ou python	Pascal Amsili, Pascal Denis	Facile à très difficile	2 personnes
		Construction automatique d'un thésaurus distributionnel du français	<p>Un thésaurus distributionnel est un dictionnaire dont chaque entrée est associée à une liste "voisins" distributionnels (à savoir des mots qui apparaissent fréquemment dans les mêmes contextes). En plus de constituer une ressource intéressante pour la linguistique de corpus, ce type de thésaurus se montre extrêmement utile (comme complément ou même comme substitut à une ressource statique telle que Wordnet) pour de nombreuses tâches de TAL (parsing syntaxique, désambiguïsation lexicale...).</p> <p>La construction d'un thésaurus distributionnel comprend trois grandes étapes: (i) l'extraction et le prétraitement d'un corpus de grande taille (plusieurs millions de mots), (ii) l'identification des contextes associés à chaque mot (ces contextes sont représentés sous la formes de simples n-grammes ou de triplets de dépendances syntaxiques), et (iii) un calcul de similarité entre contextes (ce calcul se fait typiquement sur base d'une mesure d'information mutuelle).</p> <p>Ce projet mettra donc en oeuvre cette procédure pour la construction d'un thésaurus automatique du français. En particulier, on utilisera comme corpus de départ le corpus issu de la campagne PASSAGE: un corpus de 100M de mots déjà prétraité et analysé en dépendances syntaxiques par différents analyseurs.</p> <p>Ce sujet demande des compétences solides en programmation.</p>	Java ou python	Pascal Denis	Moyen à difficile	2 personnes
		Réseau sémantique à partir d'un dictionnaire	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet. On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p>	Java ou python	Pascal Amsili	Moyen à très difficile	2 ou 3 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambigüiser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut précisément définir) des mots pleins de son contexte.</p> <p>D'autres exploitations d'un tel réseau peuvent être envisagées.</p>				
		Extraction de dépendances à partir du French Treebank	<p>Il s'agit d'explicitier les dépendances syntaxiques contenues dans les arbres syntagmatiques du French Treebank du laboratoire LLF . L'intérêt du sujet est d'une part pédagogique : il permet une familiarisation avec le traitement de corpus syntaxiquement annoté. Le passage de la constituance à la dépendance est en soi un problème linguistiquement intéressant, qui requiert une bonne maîtrise des concepts syntaxiques fondamentaux (constituants / fonctions / tête). D'autre part, le résultat de la conversion (le corpus converti en dépendances), directement exploitable pour par exemple extraire un lexique syntaxique, ou pour des expériences de parsing probabiliste en dépendances.</p> <p>La mise en œuvre consistera en :</p> <ul style="list-style-type: none"> • Etudier le corpus syntaxiquement annoté (avec par exemple l'outil tregexp). Un programme de chargement du corpus en structures python vous sera fourni. • Dégager et appliquer des règles de propagation de tête dans les constituants • Convertir les constituants en dépendances, si possible dans un format standard, comme celui de la campagne d'évaluation CoNLL 2006. <p>Vous procéderez itérativement en étudiant les erreurs / échecs à chaque itération.</p>	Python	Marie Candito	Moyen	2 personnes
2008-2009	L3	ELIZA en français	<p>On vous propose d'implémenter une version française de ELIZA. ELIZA est historiquement le premier succès spectaculaire du TAL. Il s'agit d'un petit agent conversationnel qui simule une conversation entre vous et un psychanalyste de patience infinie. Exemple de conversation en français :</p> <p>Système : "Comment allez vous ?"</p> <p>Utilisateur : "Je suis déprimé."</p> <p>S : "Dites-moi pourquoi vous êtes déprimé."</p> <p>U : "Je passe trop de temps avec les ordinateurs."</p> <p>S : "Pourquoi mentionnez-vous les ordinateurs ?"</p> <p>U : "Parce que tu en es un."</p> <p>S : "Quels sont vos rapports avec les ordinateurs ?"</p> <p>U : "Difficiles."</p> <p>S : etc</p> <p>Ce petit agent fonctionne principalement à l'aide de techniques élémentaires de pattern matching (expressions régulières) et vous demande de manipuler une pile simulant sa mémoire. Le sujet se focalisera sur les points suivants :</p> <p>Identifier un mot clé et axer la réponse sur ce mot clé (exemple : retenir un groupe de deux mots commençant par un déterminant, "les ordinateurs" dans "Je passe trop de temps avec les ordinateurs"),</p> <p>Revenir à un mot clé quand l'énoncé courant ne permet de générer aucune réponse,</p> <p>Exploiter des phrases à trou (exemple : "Pourquoi mentionnez-vous" + mot clé),</p> <p>Traduire les mots référant aux interlocuteurs : "je", "vous", "moi", "me", etc. (exemple : "je suis déprimé" devient "vous êtes déprimé" dans "Dites-moi pourquoi vous êtes déprimé").</p> <p>Sortir des réponses types au hasard (exemple : "Comment allez-vous ?").</p>	Java	Benoît Crabbé	Difficile	2 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
		Détection de la langue d'un texte	<p>Pour détecter la langue d'un texte, on peut constituer une base de connaissances à partir d'un corpus de textes classés par langue. Pour chaque langue, un premier programme (à écrire) recueillera des statistiques significatives, basées sur les lettres (par exemple, il y a plus de "w" en anglais qu'en français). Le choix du modèle probabiliste employé et de ses paramètres (bigrammes, trigrammes) devra être justifié dans le rapport. On peut en proposer plusieurs et discuter de leurs avantages et inconvénients (rapport entre précision de la reconnaissance et volume de la base de connaissances ou longueur du texte nécessaire pour reconnaître sa langue).</p> <p>Un deuxième programme (à écrire), utilisera ces bases de connaissances pour reconnaître la langue d'un texte.</p> <p>Vous utiliserez un corpus d'apprentissage comprenant des textes plus ou moins variés d'un certain nombre de langues (le plus de langues possible, au moins 4). Le corpus de test ne devra pas contenir de texte appartenant au corpus d'apprentissage.</p> <p>Amélioration possible : gérer des textes dans différents encodages, pour les langues à alphabet non latin.</p>	Java	Marie Candito	Facile à moyen	2 personnes
		Correcteur orthographique (approche graphémique)	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> de choisir un de ces remplaçants pour l'occurrence, ou pour toutes les occurrences de fournir un remplaçant pour l'occurrence, ou pour toutes les occurrences si le remplaçant n'est pas dans le lexique, de l'insérer dans le lexique pour ce texte, de l'insérer dans le lexique stable d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. <p>Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif.</p> <p>Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot :</p> <ul style="list-style-type: none"> repérage des bigrammes ou trigrammes impossibles en français redoublement ou dé-doublement de consonnes suppression/insertion de diacritiques distance d'édition etc <p>Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.</p>	Java	Marie Candito		2 personnes
		Réaccentuation	<p>Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre.</p> <p>On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple :</p> <p>La ou le français n'est pas accentué,</p> <p>il y a de la gene,</p> <p>mais quand le système m'accentue,</p>	Java	Benoît Crabbé		2 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			je suis moins gene!				
		Extraction de sigles	L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes) (par exemple EDF), et de leur forme développée (par exemple Electricité de France), utilisés dans un texte donné. Le programme aura deux fonctions : Création (extraction) de la liste (triée, normalisée) des sigles du texte donné. Attention aux formes diverses d'un même sigle (Inalf, I.N.A.L.F, INaLF). Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte. Le texte fourni est en ASCII brut	Java	Marie Candito		2 personnes
		Justification et césure	Il s'agit de réaliser un programme qui justifie (au sens des traitements de texte) un texte fourni en ASCII, sur un nombre de colonnes donné, en ajoutant des espaces entre les mots et/ou en découpant les mots selon les règles en usage pour le français. On supposera que tous les caractères ont la même dimension (« fonte fixe »). Pour le découpage (éventuel) des mots (césure), on prendra bien garde de distinguer les règles, qui seront stockées dans un fichier, du programme lui-même. On s'autorisera à utiliser un dictionnaire d'exceptions.	Java	Benoît Crabbé		2 personnes
	M1	Concordancier et outils statistiques pour corpus annoté	On propose de réaliser une suite d'outils pour l'exploration de corpus annotés en morphosyntaxe. Celle-ci se composera de deux sous-composantes : Une suite d'outils statistiques qui permettront d'obtenir des informations quantitatives sur le corpus, comme la fréquence et le nombre d'occurrences des mots. Un outil statistique qui proposera de détecter les collocations en corpus en utilisant une heuristique guidée par l'information mutuelle. Les outils statistiques devront produire une sortie texte qui permettra l'usage de leurs résultats dans un logiciel approprié au traitement de données comme Microsoft Excell ou R. Un concordancier. Le concordancier permettra de faire des recherches sur le corpus annoté. Celui-ci sera réalisé en prenant en compte trois facteurs principaux (1) L'expressivité du langage de requête (2) L'efficacité de la recherche en terme de temps de réponse et (3) de fournir une interface utilisateur conviviale permettant notamment à celui-ci de trier les résultats d'une concordance. On veillera à ce que le concordancier soit relativement indépendant des formats de corpus sur lequel il sera développé.	Java	Benoît Crabbé	Très facile à très difficile	2 personnes
		Extraction de collocations en corpus	Les collocations sont des ensembles de deux ou plusieurs mots qui sont plus fréquemment cooccurrents que la normale. Il s'agit de constructions qui ont une syntaxe classique mais une sémantique qui n'est pas complètement compositionnelle. Par exemple on considère que thé fort est une collocation car ces deux mots apparaissent assez souvent en même temps et le sens de fort dans cette séquence est légèrement altéré : il ne s'agit pas d'une grande force physique mais plutôt d'une forte concentration d'un agent du thé. Ce sujet propose de partir à la découverte de différentes méthodes connues de détection de collocations en corpus reposant sur des tests d'hypothèses statistiques (comme par exemple le test du Chi2 ou le test exact de Fisher) en utilisant des logiciels d'analyse de données appropriés.	Python ou perl et initiation à R	Benoît Crabbé	Moyen	2 personnes
		Repérage des entités nommées	On regroupe sous le terme "entités nommées" les noms de personnes, de lieux, de dates, noms d'entreprises, adresses, etc. Il s'agit d'expressions qui dénotent une entité unique de façon presque indépendante du contexte. On s'intéresse aux entités nommées pour plusieurs raisons : elles constituent des syntagmes qui peuvent être relativement complexes au point de vue syntaxique (par exemple une adresse, ou un nom d'association) dont le repérage préalable peut grandement simplifier une analyse syntaxique ; dans une perspective de recherche d'information, la reconnaissance des entités nommées permet de savoir de quoi parle un texte ; elles sont nécessaires pour la résolution des anaphores.	Java	Marie Candito	Facile à très difficile	2 ou 3 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>Il s'agit dans ce projet de repérer de la façon la plus complète possible dans un texte étiqueté ou non, les entités de type "personne". Pour cela, on envisagera un algorithme en deux étapes (qui peuvent se répéter) :</p> <p>au moyen de règles générales et de dictionnaires spécialisés (noms propres, amorces --- c'est-à-dire mots qui introduisent systématiquement des entités nommées, comme 'Melle', etc.), constitution d'une "table des symboles" des entités présentes dans le texte ;</p> <p>à partir de cette table des symboles, et en tenant compte des formes variées sous lesquelles une même entités peut être désignée, recherche de nouvelles entités, voire de nouvelles règles trouvées précédemment.</p> <p>L'idée est que le programme s'enrichit au fur et à mesure qu'il est utilisé.</p>				
		Réseaux sémantiques	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet. On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p> <p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambiguïser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut précisément définir) des mots pleins de son contexte.</p> <p>D'autres exploitations d'un tel réseau peuvent être envisagées.</p>	Java	Pascal Amsili	Moyen à très difficile	2 ou 3 personnes
		La sous-catégorisation verbale en corpus	<p>Il s'agit de compléter le French treebank du laboratoire LLF. voir http://www.llf.cnrs.fr/Gens/Abeille/French-Treebank-fr.php</p> <p>A partir de Treelex, un lexique de valence extrait par A. Kupsc à partir de ce corpus (voir http://erssab.u-bordeaux3.fr/article.php3?id_article=150) il s'agit de marquer les attributs subcat pour chaque occurrence verbale dans le corpus en désambiguïsant en contexte si plusieurs valences sont disponibles. Il s'agira d'un marquage automatique (langage de programmation libre, python recommandé) à valider manuellement. Il s'agira aussi de marquer la voix passive le cas échéant.</p> <p>Le projet peut se concerner au marquage automatique, et la validation manuelle peut se faire dans le cadre d'un stage rémunéré par le laboratoire LLF (300 euros net/mois). Il est validé deux fois.</p>		Anne Abeillé, Marie Candito	Facile à moyen	1 personne
		Marquage des antécédents pronominaux en corpus	<p>Il s'agit de tester des algorithmes de résolution d'anaphores sur texte analysé syntaxiquement (Le French treebank du laboratoire LLF pourra être utilisé, voir http://www.llf.cnrs.fr/Gens/Abeille/French-Treebank-fr.php</p> <p>A partir des occurrences pronominales de 3e personne marquées pour le genre et le nombre de leurs antécédents, il s'agit de proposer un ou plusieurs candidats dans la même phrase (pour les relatifs et les réfléchis) ou non (pour les personnels).</p> <p>Les mots antécédents seront repérés par leur numéro dans la phrase (ID) éventuellement préfixé par le numéro de la phrase.</p>	Java ou python	Pascal Amsili	Facile à très difficile	2 personnes
2007-2008	L3						
	M1						
2006-2007	L3						
	M1						
2005-2006	L3						

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
2000	M1	Repérage des entités nommées	<p>On regroupe sous le terme "entités nommées" les noms de personnes, de lieux, de dates, noms d'entreprises, adresses, etc. Il s'agit d'expressions qui dénotent une entité unique de façon presque indépendante du contexte. On s'intéresse aux entités nommées pour plusieurs raisons :</p> <ul style="list-style-type: none"> o elles constituent des syntagmes qui peuvent être relativement complexes au point de vue syntaxique (par exemple une adresse, ou un nom d'association) dont le repérage préalable peut grandement simplifier une analyse syntaxique ; o dans une perspective de recherche d'information, la reconnaissance des entités nommées permet de savoir de quoi parle un texte ; o elles sont nécessaires pour la résolution des anaphores. <p>Il s'agit dans ce projet de repérer de la façon la plus complète possible dans un texte étiqueté ou non, les entités temporelles. Pour cela, on envisagera un algorithme en deux étapes (qui peuvent se répéter) :</p> <ul style="list-style-type: none"> o au moyen de règles générales et de dictionnaires spécialisés (noms propres, amorces --- c'est-à-dire mots qui introduisent ou suivent systématiquement des entités nommées, comme 'minutes', etc.), constitution d'une "table des symboles" des entités présentes dans le texte ; o à partir de cette table des symboles, et en tenant compte des formes variées sous lesquelles une même entité peut être désignée, recherche de nouvelles entités, voire de nouvelles règles trouvées précédemment. <p>L'idée est que le programme s'enrichit au fur et à mesure qu'il est utilisé.</p>	<i>Au choix parmi les langages étudiés en cours</i>			
		Analyseur syntaxique d'une grammaire ambiguë (Earley)	Il s'agit d'implémenter, pour un fragment significatif du français (ou d'une autre langue), comportant des ambiguïtés, l'algorithme d'analyse d'Earley.	<i>Au choix parmi les langages étudiés en cours</i>			
		Réseaux sémantiques	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet.</p> <p>On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p> <p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambiguïser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut précisément définir) des mots pleins de son contexte.</p> <p>D'autres exploitations d'un tel réseau peuvent être envisagées.</p>	<i>Au choix parmi les langages étudiés en cours</i>			
		Détection de la langue d'un texte	<p>Pour détecter la langue d'un texte, on peut constituer une base de connaissances à partir d'un corpus de textes classés par langue. Pour chaque langue, un premier programme (à écrire) recueillera des statistiques significatives, basées sur les lettres (par exemple, il y a plus de "w" en anglais qu'en français). Le choix du modèle probabiliste employé et de ses paramètres (bigrammes, trigrammes) devra être justifié dans le rapport. On peut en proposer plusieurs et discuter de leurs avantages et inconvénients (rapport entre précision de la reconnaissance et volume de la base de connaissances ou longueur du texte nécessaire pour reconnaître sa langue).</p> <p>Un deuxième programme (à écrire), utilisera ces bases de connaissances pour reconnaître la langue d'un texte.</p>	<i>Au choix parmi les langages étudiés en cours</i>			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>Vous utiliserez un corpus d'apprentissage comprenant des textes plus ou moins variés d'un certain nombre de langues (le plus de langues possible, au moins 4). Le corpus de test ne devra pas contenir de texte appartenant au corpus d'apprentissage.</p> <p>Amélioration possible : gérer des textes dans différents encodages, pour les langues à alphabet non latin.</p>				
		Wordnet en java		Java			
		Résolution des anaphores personnelles		<i>Au choix parmi les langages étudiés en cours</i>			
2004-2005	L3	Correcteur orthographique (approche graphémique)	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> • de choisir un de ces remplaçants <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences • de fournir un remplaçant <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences o si le remplaçant n'est pas dans le lexique, <input type="checkbox"/> de l'insérer dans le lexique pour ce texte, <input type="checkbox"/> de l'insérer dans le lexique stable • d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. <p>Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif.</p> <p>Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot :</p> <ul style="list-style-type: none"> • repérage des bigrammes ou trigrammes impossibles en français • redoublement ou dé-doublement de consonnes • suppression/insertion de diacritiques • distance d'édition • etc. <p>Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.</p>	Pascal ou C			
		Réaccentuation	<p>Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre.</p> <p>On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple :</p> <p>La ou le francais n'est pas accentue, il y a de la gene, mais quand le systeme m'accentue, je suis moins gene!</p>	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
		Extraction de sigles	<p>L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes) (par exemple EDF), et de leur forme développée (par exemple Electricité de France), utilisés dans un texte donné. Le programme aura deux fonctions :</p> <ol style="list-style-type: none"> 1. Création (extraction) de la liste (triée, normalisée) des sigles du texte donné. Attention aux formes diverses d'un même sigle (Inalf, I.N.A.L.F, INaLF). 2. Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte. <p>Le texte fourni est en ASCII brut.</p>	Pascal ou C			
		Justification et césure	<p>Il s'agit de réaliser un programme qui justifie (au sens des traitements de texte) un texte fourni en ASCII, sur un nombre de colonnes donné, en ajoutant des espaces entre les mots et/ou en découpant les mots selon les règles en usage pour le français.</p> <p>On supposera que tous les caractères ont la même dimension (« fonte fixe »).</p> <p>Pour le découpage (éventuel) des mots (césure), on prendra bien garde de distinguer les règles, qui seront stockées dans un fichier, du programme lui-même. On s'autorisera à utiliser un dictionnaire d'exceptions.</p>	Pascal ou C			
		Marquage des syntagmes nominaux	<p>Il s'agit, dans un texte étiqueté, de repérer le début et la fin des syntagmes nominaux (ou d'une bonne partie d'entre eux), de manière à produire en sortie un texte dans lequel ces débuts et fin de SN seront balisés (par exemple <sn>le chat de <sn>la voisine</sn></sn>).</p> <p>Le programme permettra à la fois une utilisation "batch", qui créera directement le fichier résultat, et une utilisation interactive, qui permettra à l'utilisateur de voir les résultats produits au fur et à mesure.</p> <p>Une étude linguistique va permettre de déterminer les séquences de lemmes ou d'étiquettes qui vont correspondre à un début ou une fin de syntagme. On fera en sorte que ces "règles" soient placées dans un fichier à part, et non "codées en dur", pour permettre une mise au point plus facile et une meilleure qualité du programme.</p>	Pascal ou C			
	M1	Repérage des entités nommées	<p>On regroupe sous le terme "entités nommées" les noms de personnes, de lieux, de dates, noms d'entreprises, adresses, etc. Il s'agit d'expressions qui dénotent une entité unique de façon presque indépendante du contexte. On s'intéresse aux entités nommées pour plusieurs raisons :</p> <ul style="list-style-type: none"> • elles constituent des syntagmes qui peuvent être relativement complexes au point de vue syntaxique (par exemple une adresse, ou un nom d'association) dont le repérage préalable peut grandement simplifier une analyse syntaxique ; • dans une perspective de recherche d'information, la reconnaissance des entités nommées permet de savoir de quoi parle un texte ; • elles sont nécessaires pour la résolution des anaphores. <p>Il s'agit dans ce projet de repérer de la façon la plus complète possible dans un texte étiqueté ou non, les entités de type "personne". Pour cela, on envisagera un algorithme en deux étapes (qui peuvent se répéter) :</p> <ul style="list-style-type: none"> • au moyen de règles générales et de dictionnaires spécialisés (noms propres, amorces --- c'est-à-dire mots qui introduisent systématiquement des entités nommées, comme 'Melle', etc.), constitution d'une "table des symboles" des entités présentes dans le texte ; • à partir de cette table des symboles, et en tenant compte des formes variées sous lesquelles une même entités peut être désignée, recherche de nouvelles entités, voire de nouvelles règles trouvées précédemment. <p>L'idée est que le programme s'enrichit au fur et à mesure qu'il est utilisé.</p>	Lisp, Prolog ou java			
		Analyseur syntaxique d'une grammaire ambiguë (Earley)	<p>Il s'agit d'implémenter, pour un fragment significatif du français (ou d'une autre langue), comportant des ambiguïtés, l'algorithme d'analyse d'Earley.</p>	Lisp, Prolog ou java			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
		Réseaux sémantiques	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet.</p> <p>On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p> <p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambigüiser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut préciser) des mots pleins de son contexte.</p> <p>D'autres exploitations d'un tel réseau peuvent être envisagées.</p>	Lisp, Prolog ou java			
		Détection de la langue d'un texte	<p>Pour détecter la langue d'un texte, on peut constituer une base de connaissances à partir d'un corpus de textes classés par langue. Pour chaque langue, un premier programme (à écrire) recueillera des statistiques significatives, basées sur les lettres (par exemple, il y a plus de "w" en anglais qu'en français). Le choix du modèle probabiliste employé et de ses paramètres (bigrammes, trigrammes) devra être justifié dans le rapport. On peut en proposer plusieurs et discuter de leurs avantages et inconvénients (rapport entre précision de la reconnaissance et volume de la base de connaissances ou longueur du texte nécessaire pour reconnaître sa langue).</p> <p>Un deuxième programme (à écrire), utilisera ces bases de connaissances pour reconnaître la langue d'un texte.</p> <p>Vous utiliserez un corpus d'apprentissage comprenant des textes plus ou moins variés d'un certain nombre de langues (le plus de langues possible, au moins 4). Le corpus de test ne devra pas contenir de texte appartenant au corpus d'apprentissage.</p> <p>Amélioration possible : gérer des textes dans différents encodages, pour les langues à alphabet non latin.</p>	Lisp, Prolog ou java			
		Wordnet en java		Java			
2003-2004	L3						
	M1						
2002-2003	L3	Correcteur orthographique (approche graphémique)	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> • de choisir un de ces remplaçants <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences • de fournir un remplaçant <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences o si le remplaçant n'est pas dans le lexique, <input type="checkbox"/> de l'insérer dans le lexique pour ce texte, <input type="checkbox"/> de l'insérer dans le lexique stable • d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. 	Pascal ou C			2 ou 3 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif.</p> <p>Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot :</p> <ul style="list-style-type: none"> • repérage des bigrammes ou trigrammes impossibles en français • redoublement ou dé-doublement de consonnes • suppression/insertion de diacritiques • distance d'édition • etc. <p>Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.</p>				
		Réaccentuation	<p>Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre.</p> <p>On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple :</p> <p>La ou le francais n'est pas accentuee, il y a de la gene, mais quand le systeme m'accentuee, je suis moins gene!</p> <p>On pourra s'inspirer du programme « Reacc » développé par le RALI à l'Université de Montréal.</p>	Pascal ou C			2 ou 3 personnes
		Recherche de patrons	<p>L'objectif de ce projet est de réaliser un programme permettant de détecter dans un corpus étiqueté et lemmatisé des séquences de mots conformes à des patrons lexico-syntaxiques. Un patron étant une suite de mots, de lemmes ou de catégories morpho-syntaxiques. Exemples: « N à N », « relation de N »</p> <p>Le programme prendra en entrée un corpus et un fichier de patrons et produira en sortie une liste classée par fréquence, de toutes les séquences de mots conformes aux patrons, accompagnées de leur nombre d'occurrences.</p> <p>Dans la définition des patrons, il faudra introduire le moyen de distinguer les mots des lemmes et des catégories.</p> <p>On attachera un soin particulier au choix de la structure de données, et à l'algorithme, qui doit faire le calcul en un temps raisonnable.</p>	Pascal ou C			2 ou 3 personnes
		Extraction de sigles	<p>L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes) (par exemple EDF), et de leur forme développée (par exemple Electricité de France), utilisés dans un texte donné. Le programme aura deux fonctions :</p> <ol style="list-style-type: none"> 1. Création (extraction) de la liste (triée, normalisée) des sigles du texte donné. Attention aux formes diverses d'un même sigle (Inalf, I.N.A.L.F, INaLF). 2. Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte. <p>Le texte fourni est en ASCII brut.</p>	Pascal ou C			2 ou 3 personnes
		Justification et césure	<p>Il s'agit de réaliser un programme qui justifie (au sens des traitements de texte) un texte fourni en ASCII, sur un nombre de colonnes donné, en ajoutant des espaces entre les mots et/ou en découpant les mots selon les règles en usage pour le français.</p> <p>On supposera que tous les caractères ont la même dimension (« fonte fixe »).</p> <p>Pour le découpage (éventuel) des mots (césure), on prendra bien garde de distinguer les règles, qui seront stockées dans un fichier, du programme lui-même. On s'autorisera à utiliser un dictionnaire d'exceptions.</p>	Pascal ou C			2 ou 3 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
		Enrichissement de requêtes	<p>Il s'agit de se placer dans la problématique suivante : on formule des requêtes à un moteur de recherche, les requêtes pouvant contenir un ou plusieurs mots, voire un syntagme ou une phrase. Sachant que les taux de réussite de telles requêtes ne sont pas toujours satisfaisants, on se propose d'enrichir ces requêtes (en ajoutant des mots, ou en utilisant les opérateurs classiques de moteurs de recherche...) pour améliorer le résultat.</p> <p>Cet enrichissement sera réalisé sur la bases d'information de nature lexicale ou terminologique, stockées dans un réseau sémantique approprié.</p> <p>Le but du projet n'est pas de réaliser un tel réseau sémantique à grande couverture, mais de spécifier, pour un domaine conceptuel fixé, les heuristiques à utiliser et l'architecture du réseau correspondant. Le projet donnera lieu à la réalisation d'une maquette, au sens industriel du terme.</p>	Pascal ou C			2 ou 3 personnes
	M1	Repérage des "entités nommées"	<p>On regroupe sous le terme « entités nommées » les noms de personnes, de lieux, les dates, noms d'entreprises, adresses, etc. Le but du projet consiste à repérer automatiquement un sous-ensemble homogène parmi les entités (par exemple les noms de personnes ou les noms d'entreprises), l'idée étant de faciliter un traitement ultérieur, sémantique (extraction d'information, résolution d'anaphore...), syntaxique (analyse, chunking...), etc.</p> <p>Le programme envisagé procède par étapes successives sur le même texte. Dans une première étape, grâce aux dictionnaires spécialisés (cf. plus loin), on étiquette les mots simples. On partira d'un corpus déjà étiqueté morpho-syntaxiquement. La seconde étape, basée sur des règles, permettra de regrouper ces mots simples pour former des entités nommées. Enfin, une troisième étape, facultative dans ce projet, pourrait tenter de relier entre elles les entités co-référentes. Les éléments seront marqués dans le texte au moyen d'un système de balisage de type SGML.</p> <p>Le système repose donc, outre le programme lui-même, sur les ressources suivantes :</p> <ul style="list-style-type: none"> • des dictionnaires : <ul style="list-style-type: none"> o dictionnaire de la langue générale pour repérer les mots inconnus. o dictionnaires spécialisés de prénoms, de noms de sociétés connues, etc. o dictionnaires spécialisés d'« amorces », mots ou groupes de mots qui peuvent marquer le début d'une entité nommée. • une « grammaire » permettant de combiner des mots repérés de natures diverses (mots inconnus, prénoms, amorces...) pour en faire une entité (un nom de personne). <p>Les dictionnaires et la grammaire sont stockés dans des fichiers textes. On assure ainsi une séparation nette entre les données et les traitement (le programme).</p> <p>Exemple</p> <p>Soit le texte « M. Jack Lang a gagné la mairie. »</p> <p>"M." est stocké dans un dictionnaire "d'amorces"</p> <p>"Jack" est stocké dans un dictionnaire de prénoms</p> <p>"Lang" est inconnu du système</p> <p>La première étape permet d'obtenir :</p> <pre><TR_PERSON>M.</TR_PERSON> <PERSON>Jack</PERSON> <UFIRSTUNKNOWN>Lang</UFIRSTUNKNOWN></pre> <p>a gagné la mairie.</p> <p>Une règle de la forme</p> <pre>TR_PERSON PERSON UFIRSTUNKNOWN ==> PERSON</pre> <p>va pouvoir être activée.</p> <p>D'où le résultat de la seconde étape :</p> <pre><PERSON>M. Jack Lang</PERSON> a gagné la mairie.</pre> <p>Une troisième étape permettrait d'identifier toutes les occurrences de Jack Lang</p>	Lisp, Prolog ou java			2 ou 3 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<PERSON id=1>M. Jack Lang</PERSON> a gagné la mairie.				
		Réseaux sémantiques	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet.</p> <p>On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p> <p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambiguïser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut précisément définir) des mots pleins de son contexte.</p> <p>D'autres exploitations d'un tel réseau peuvent être envisagées.</p>	Lisp, Prolog ou java			2 ou 3 personnes
		"Complétion" automatique	<p>On se propose de réaliser un programme qui accélère la vitesse de frappe d'un utilisateur en proposant - de façon intelligente - la fin des mots en cours de frappe.</p> <p>Un tel système, pour être vraiment performant, devrait faire de la prédiction syntaxique en temps réel. On se contentera dans ce projet de faire un système qui propose la suite la plus probable des mots.</p> <p>La méthode repose sur un calcul de probabilité basé sur la fréquence d'apparition de chaque mot dans un corpus préalablement fourni. On construit un arbre dont chaque chemin correspond à un mot (une lettre par arc), et dont chaque arc porte une probabilité.</p> <p>Par exemple, la probabilité de la prochaine lettre de la suite ``v-o-l-u" est respectivement de 0/9, 8/9 (3+5/9), 1/9, 0/9 pour les lettres `b', `m', `p' et `t' si nous avons dans le corpus :</p> <p>0 occurrences du mot volubile 3 " volume 5 " volumineux 1 " voluptueux 0 " volute</p> <p>La lettre la plus probable, ici ``m", sera proposée à l'utilisateur, et un nouveau calcul est effectué pour la suite de lettres ``v-o-l-u-m".</p> <p>Un soin particulier sera apporté à l'interface de ce programme : l'utilisateur doit pouvoir accepter ou revenir facilement sur les propositions du système, etc. La dimension ergonomique de ce projet est au moins aussi importante que la dimension algorithmique : en particulier, il conviendra de choisir entre la stratégie qui consiste à ne produire automatiquement qu'une lettre à la fois, auquel cas notre objectif d'économie en vitesse de frappe est compromis et celle qui consiste à produire le suffixe en entier (à partir de quelle longueur de préfixe ?), au risque de nécessiter beaucoup de corrections de la part de l'utilisateur.</p> <p>Un langage particulièrement adapté au projet serait Emacs Lisp, ce qui permettrait d'intégrer la complétion à l'éditeur Emacs, ce qui n'exclut cependant pas la réalisation de ce projet dans un autre langage.</p>	Emacs Lisp			2 ou 3 personnes
		Étiqueteur morpho-syntaxique basé sur un tri-gramme	<p>L'étiqueteur morpho-syntaxique développé dans le cadre du cours LI031 repose sur un modèle bi-gramme : la prédiction de la catégorie d'un mot dépend des catégories possibles du mot et de la catégorie du mot précédent.</p> <p>L'objectif de ce projet est de modifier l'étiqueteur pour que la prédiction repose sur la catégorie des deux mots précédents.</p>	Lisp, Prolog ou java			2 ou 3 personnes

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			Les modifications porteront d'une part sur les structures de données permettant de stocker les tri-grammes de façon efficace et d'autre part sur l'implémentation des méthodes de repli et d'interpolation permettant de pallier l'insuffisance des données d'apprentissage, plus importantes que dans le cas d'un bi-gramme.				
		Traitement des mots inconnus dans un étiqueteur morpho-syntaxique	<p>L'étiqueteur morpho-syntaxique développé dans le cadre du cours LI031 ne prévoit pas le cas où un mot d'une phrase à étiqueter ne figure pas dans le lexique de l'étiqueteur. L'objectif de ce projet est d'intégrer à l'étiqueteur un module de gestion des mots inconnus.</p> <p>Plusieurs solutions seront testées, parmi lesquelles, la prédiction de la catégorie du mot inconnu en fonction de sa forme (en particulier, les préfixes et les suffixes) et en fonction de son contexte (les catégories possibles des mots suivants et précédents) ou encore de la loi de probabilité qu'un mot inconnu appartienne aux différentes catégories (un mot inconnu à une plus grande probabilité d'être un nom pronom qu'un article.). Les différentes solutions seront testées isolément et conjointement, une analyse des erreurs permettra de déterminer les points forts et les points faibles des différentes solutions.</p>	Lisp, Prolog ou java			2 ou 3 personnes
		Désambiguïsation par automates	<p>Partant d'un texte dont chaque mot est étiqueté avec toutes ses catégories morpho-syntaxiques possibles (étiquetage sans désambiguïsation), on va s'attacher à lever certaines ambiguïtés non pertinentes, c'est-à-dire à supprimer certaines des étiquettes attachées à certains mots.</p> <p>L'idée est de disposer d'un programme capable de prendre en entrée des règles comme la suivante : « si X est marqué N et V, et que X est précédé d'un Y marqué Det, alors l'étiquette V doit être ôtée à X » ; et d'appliquer ces règles sur la totalité du corpus.</p> <p>On conseille d'implémenter ces règles sous forme d'automates.</p> <p>Le travail comprend donc une partie linguistique, visant à mettre en évidence des règles de désambiguïsation pour le français, basées sur les étiquettes et/ou les lemmes des mots précédents et suivants les mots ambigus une partie formelle, qui s'attachera à établir un format de règles correspondant aux besoins relevés par l'étude linguistique et une partie informatique, qui permettra à partir d'un corpus quelconque, et d'un fichier de règles, de désambiguïser le corpus (ou du moins d'en diminuer l'ambiguïté).</p>	Lisp, Prolog ou java			2 ou 3 personnes
2001-2002	L3	Correcteur orthographique (approche graphémique)	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> • de choisir un de ces remplaçants <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences • de fournir un remplaçant <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences o si le remplaçant n'est pas dans le lexique, <input type="checkbox"/> de l'insérer dans le lexique pour ce texte, <input type="checkbox"/> de l'insérer dans le lexique stable • d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. <p>Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif.</p> <p>Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot :</p> <ul style="list-style-type: none"> • repérage des bigrammes ou trigrammes impossibles en français • redoublement ou dé-doublement de consonnes • suppression/insertion de diacritiques • distance d'édition 	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<ul style="list-style-type: none"> • etc. <p>Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.</p>				
		Réaccentuation	<p>Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre.</p> <p>On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple :</p> <p>La ou le francais n'est pas accentue, il y a de la gene, mais quand le systeme m'accentue, je suis moins gene!</p> <p>On pourra s'inspirer du programme « Reacc » développé par le RALI à l'Université de Montréal.</p>	Pascal ou C			
		Recherche de patrons	<p>L'objectif de ce projet est de réaliser un programme permettant de détecter dans un corpus étiqueté et lemmatisé des séquences de mots conformes à des patrons lexico-syntaxiques. Un patron étant une suite de mots, de lemmes ou de catégories morpho-syntaxiques. Exemples: « N à N », « relation de N »</p> <p>Le programme prendra en entrée un corpus et un fichier de patrons et produira en sortie une liste classée par fréquence, de toutes les séquences de mots conformes aux patrons, accompagnées de leur nombre d'occurrences.</p> <p>Dans la définition des patrons, il faudra introduire le moyen de distinguer les mots des lemmes et des catégories.</p> <p>On attachera un soin particulier au choix de la structure de données, et à l'algorithme, qui doit faire le calcul en un temps raisonnable.</p>	Pascal ou C			
		Extraction de sigles	<p>L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes) (par exemple EDF), et de leur forme développée (par exemple Electricité de France), utilisés dans un texte donné. Le programme aura deux fonctions :</p> <ol style="list-style-type: none"> 1. Création (extraction) de la liste (triée, normalisée) des sigles du texte donné. Attention aux formes diverses d'un même sigle (Inalf, I.N.A.L.F, INaLF). 2. Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte. <p>Le texte fourni est en ASCII brut.</p>	Pascal ou C			
		Justification et césure	<p>Il s'agit de réaliser un programme qui justifie (au sens des traitements de texte) un texte fourni en ASCII, sur un nombre de colonnes donné, en ajoutant des espaces entre les mots et/ou en découpant les mots selon les règles en usage pour le français.</p> <p>On supposera que tous les caractères ont la même dimension (« fonte fixe »).</p> <p>Pour le découpage (éventuel) des mots (césure), on prendra bien garde de distinguer les règles, qui seront stockées dans un fichier, du programme lui-même. On s'autorisera à utiliser un dictionnaire d'exceptions.</p>	Pascal ou C			
		Enrichissement de requêtes	<p>Il s'agit de se placer dans la problématique suivante : on formule des requêtes à un moteur de recherche, les requêtes pouvant contenir un ou plusieurs mots, voire un syntagme ou une phrase. Sachant que les taux de réussite de telles requêtes ne sont pas toujours satisfaisants, on se propose d'enrichir ces requêtes (en ajoutant des mots, ou en utilisant les opérateurs classiques de moteurs de recherche...) pour améliorer le résultat.</p> <p>Cet enrichissement sera réalisé sur la bases d'information de nature lexicale ou terminologique, stockées dans un réseau sémantique approprié.</p>	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			Le but du projet n'est pas de réaliser un tel réseau sémantique à grande couverture, mais de spécifier, pour un domaine conceptuel fixé, les heuristiques à utiliser et l'architecture du réseau correspondant. Le projet donnera lieu à la réalisation d'une maquette, au sens industriel du terme.				
		Téléphone	<p>La plupart des téléphones associent à chaque touche un certain nombre de lettres, permettant ainsi de transmettre des messages. Du fait qu'un chiffre ne correspond pas à une seule lettre, une suite de chiffres peut être ambiguë (elle peut correspondre à plus d'un mot). La suite 7-6-8-7, par exemple, correspond aux mots pour, sous et pots.</p> <p>Le but du projet est de réaliser un programme qui prend en entrée une suite de chiffre et un dictionnaire de formes fléchies et qui propose en sortie la liste des mots du dictionnaire qui correspondent à la suite des chiffres. Ces mots seront classés selon une fréquence représentée dans le dictionnaire.</p> <p>On attachera un soin particulier au choix de la structure de données, et à l'algorithme, qui doit faire le calcul en un temps raisonnable.</p>	Pascal ou C			
	M1	Repérage des "entités nommées"	<p>On regroupe sous le terme « entités nommées » les noms de personnes, de lieux, les dates, noms d'entreprises, adresses, etc. Le but du projet consiste à repérer automatiquement un sous-ensemble homogène parmi les entités (par exemple les noms de personnes ou les noms d'entreprises), l'idée étant de faciliter un traitement ultérieur, sémantique (extraction d'information, résolution d'anaphore...), syntaxique (analyse, chunking...), etc.</p> <p>Le programme envisagé procède par étapes successives sur le même texte. Dans une première étape, grâce aux dictionnaires spécialisés (cf. plus loin), on étiquette les mots simples. On partira d'un corpus déjà étiqueté morpho-syntaxiquement. La seconde étape, basée sur des règles, permettra de regrouper ces mots simples pour former des entités nommées. Enfin, une troisième étape, facultative dans ce projet, pourrait tenter de relier entre elles les entités co-référentes. Les éléments seront marqués dans le texte au moyen d'un système de balisage de type SGML.</p> <p>Le système repose donc, outre le programme lui-même, sur les ressources suivantes :</p> <ul style="list-style-type: none"> • des dictionnaires : <ul style="list-style-type: none"> o dictionnaire de la langue générale pour repérer les mots inconnus. o dictionnaires spécialisés de prénoms, de noms de sociétés connues, etc. o dictionnaires spécialisés d'« amorces », mots ou groupes de mots qui peuvent marquer le début d'une entité nommée. • une « grammaire » permettant de combiner des mots repérés de natures diverses (mots inconnus, prénoms, amorces...) pour en faire une entité (un nom de personne). <p>Les dictionnaires et la grammaire sont stockés dans des fichiers textes. On assure ainsi une séparation nette entre les données et les traitements (le programme).</p> <p>Exemple</p> <p>Soit le texte « M. Jack Lang a gagné la mairie. »</p> <p>"M." est stocké dans un dictionnaire "d'amorces"</p> <p>"Jack" est stocké dans un dictionnaire de prénoms</p> <p>"Lang" est inconnu du système</p> <p>La première étape permet d'obtenir :</p> <pre><TR_PERSON>M.</TR_PERSON> <PERSON>Jack</PERSON> <UFIRSTUNKNOWN>Lang</UFIRSTUNKNOWN></pre> <p>a gagné la mairie.</p> <p>Une règle de la forme</p> <pre>TR_PERSON PERSON UFIRSTUNKNOWN ==> PERSON</pre> <p>va pouvoir être activée.</p> <p>D'où le résultat de la seconde étape :</p>	Lisp, Prolog, Pascal, C, C++ ou java			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p><PERSON>M. Jack Lang</PERSON> a gagné la mairie.</p> <p>Une troisième étape permettrait d'identifier toutes les occurrences de Jack Lang</p> <p><PERSON id=1>M. Jack Lang</PERSON> a gagné la mairie.</p>				
		Réseaux sémantiques	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet.</p> <p>On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p> <p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambiguïser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut précisément définir) des mots pleins de son contexte.</p> <p>D'autres exploitations d'un tel réseau peuvent être envisagées.</p>	Lisp, Prolog, Pascal, C, C++ ou java			
		"Complétion" automatique	<p>On se propose de réaliser un programme qui accélère la vitesse de frappe d'un utilisateur en proposant - de façon intelligente - la fin des mots en cours de frappe.</p> <p>Un tel système, pour être vraiment performant, devrait faire de la prédiction syntaxique en temps réel. On se contentera dans ce projet de faire un système qui propose la suite la plus probable des mots.</p> <p>La méthode repose sur un calcul de probabilité basé sur la fréquence d'apparition de chaque mot dans un corpus préalablement fourni. On construit un arbre dont chaque chemin correspond à un mot (une lettre par arc), et dont chaque arc porte une probabilité.</p> <p>Par exemple, la probabilité de la prochaine lettre de la suite ``v-o-l-u" est respectivement de 0/9, 8/9 (3+5/9), 1/9, 0/9 pour les lettres `b', `m', `p' et `t' si nous avons dans le corpus :</p> <p>0 occurrences du mot volubile</p> <p>3 " volume</p> <p>5 " volumineux</p> <p>1 " voluptueux</p> <p>0 " volute</p> <p>La lettre la plus probable, ici ``m", sera proposée à l'utilisateur, et un nouveau calcul est effectué pour la suite de lettres ``v-o-l-u-m".</p> <p>Un soin particulier sera apporté à l'interface de ce programme : l'utilisateur doit pouvoir accepter ou revenir facilement sur les propositions du système, etc. La dimension ergonomique de ce projet est au moins aussi importante que la dimension algorithmique : en particulier, il conviendra de choisir entre la stratégie qui consiste à ne produire automatiquement qu'une lettre à la fois, auquel cas notre objectif d'économie en vitesse de frappe est compromis et celle qui consiste à produire le suffixe en entier (à partir de quelle longueur de préfixe ?), au risque de nécessiter beaucoup de corrections de la part de l'utilisateur.</p> <p>Un langage particulièrement adapté au projet serait Emacs Lisp, ce qui permettrait d'intégrer la complétion à l'éditeur Emacs, ce qui n'exclut cependant pas la réalisation de ce projet dans un autre langage.</p>	Emacs Lisp			
		Étiqueteur morpho-syntaxique basé sur un tri-gramme	<p>L'étiqueteur morpho-syntaxique développé dans le cadre du cours LI031 repose sur un modèle bi-gramme : la prédiction de la catégorie d'un mot dépend des catégories possibles du mot et de la catégorie du mot précédent. L'objectif de ce projet est de modifier l'étiqueteur pour que la prédiction repose sur la catégorie des deux mots précédents.</p>	Lisp, Prolog, Pascal, C, C++ ou java			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			Les modifications porteront d'une part sur les structures de données permettant de stocker les tri-grammes de façon efficace et d'autre part sur l'implémentation des méthodes de repli et d'interpolation permettant de pallier l'insuffisance des données d'apprentissage, plus importantes que dans le cas d'un bi-gramme.				
		Traitement des mots inconnus dans un étiqueteur morpho-syntaxique	<p>L'étiqueteur morpho-syntaxique développé dans le cadre du cours LI031 ne prévoit pas le cas où un mot d'une phrase à étiqueter ne figure pas dans le lexique de l'étiqueteur. L'objectif de ce projet est d'intégrer à l'étiqueteur un module de gestion des mots inconnus.</p> <p>Plusieurs solutions seront testées, parmi lesquelles, la prédiction de la catégorie du mot inconnu en fonction de sa forme (en particulier, les préfixes et les suffixes) et en fonction de son contexte (les catégories possibles des mots suivants et précédents) ou encore de la loi de probabilité qu'un mot inconnu appartienne aux différentes catégories (un mot inconnu à une plus grande probabilité d'être un nom pronom qu'un article.). Les différentes solutions seront testées isolément et conjointement, une analyse des erreurs permettra de déterminer les points forts et les points faibles des différentes solutions.</p>	Lisp, Prolog, Pascal, C, C++ ou java			
		Désambiguïsation par automates	<p>Partant d'un texte dont chaque mot est étiqueté avec toutes ses catégories morpho-syntaxiques possibles (étiquetage sans désambiguïsation), on va s'attacher à lever certaines ambiguïtés non pertinentes, c'est-à-dire à supprimer certaines des étiquettes attachées à certains mots.</p> <p>L'idée est de disposer d'un programme capable de prendre en entrée des règles comme la suivante : « si X est marqué N et V, et que X est précédé d'un Y marqué Det, alors l'étiquette V doit être ôtée à X » ; et d'appliquer ces règles sur la totalité du corpus.</p> <p>On conseille d'implémenter ces règles sous forme d'automates.</p> <p>Le travail comprend donc une partie linguistique, visant à mettre en évidence des règles de désambiguïsation pour le français, basées sur les étiquettes et/ou les lemmes des mots précédents et suivants les mots ambigus une partie formelle, qui s'attachera à établir un format de règles correspondant aux besoins relevés par l'étude linguistique et une partie informatique, qui permettra à partir d'un corpus quelconque, et d'un fichier de règles, de désambiguïser le corpus (ou du moins d'en diminuer l'ambiguïté).</p>	Lisp, Prolog, Pascal, C, C++ ou java			
2000-2001	L3	Correcteur orthographique (approche graphémique)	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> • de choisir un de ces remplaçants <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences • de fournir un remplaçant <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences o si le remplaçant n'est pas dans le lexique, ♣ de l'insérer dans le lexique pour ce texte, ♣ de l'insérer dans le lexique stable • d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. <p>Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif.</p> <p>Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot :</p> <ul style="list-style-type: none"> • repérage des bigrammes ou trigrammes impossibles en français • redoublement ou dé-doublement de consonnes • suppression/insertion de diacritiques • insertion de lettres 	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<ul style="list-style-type: none"> • interversion de lettres • etc. <p>Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.</p>				
		Réaccentuation	<p>Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre.</p> <p>On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple :</p> <p>La ou le français n'est pas accentué, il y a de la gène, mais quand le système m'accentue, je suis moins gêné!</p> <p>On pourra s'inspirer du programme « Reacc » développé par le RALI à l'Université de Montréal.</p>	Pascal ou C			
		Calcul de collocations	<p>On veut réaliser un programme de calcul des collocations dans un corpus étiqueté et lemmatisé.</p> <p>Le but est de repérer et de compter les blocs de 2 ou 3 mots juxtaposés qui se retrouvent au moins deux fois dans le corpus. On s'intéressera aussi aux groupes « N à N », « N de N », etc., dont on calculera aussi les fréquences.</p> <p>La sortie du programme consiste donc en une liste, classée par fréquence, de toutes les collocations qui y apparaissent au moins deux fois.</p> <p>On attachera un soin particulier au choix de la structure de données, et à l'algorithme, qui doit faire le calcul en un temps raisonnable.</p>	Pascal ou C			
		Extraction de sigles	<p>L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes) (par exemple EDF), et de leur forme développée (par exemple Electricité de France), utilisés dans un texte donné. Le programme aura deux fonctions :</p> <ol style="list-style-type: none"> 1. Création (extraction) de la liste (triée, normalisée) des sigles du texte donné. Attention aux formes diverses d'un même sigle (Inalf, I.N.A.L.F, INaLF). 2. Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte. <p>Le texte fourni est en ASCII brut.</p>	Pascal ou C			
		Justification et "hyphénation"	<p>Il s'agit de réaliser un programme qui justifie (au sens des traitements de texte) un texte fourni en ASCII, sur un nombre de colonnes donné, en découpant les mots selon les règles en usage pour le français.</p> <p>On supposera que tous les caractères ont la même dimension (« fonte proportionnelle »).</p> <p>Pour le découpage (éventuel) des mots (hyphenation), on prendra bien garde de distinguer les règles, qui seront stockées dans un fichier, du programme lui-même. On s'autorisera à utiliser un dictionnaire d'exceptions.</p>	Pascal ou C			
		Enrichissement de requêtes	<p>Il s'agit de se placer dans la problématique suivante : on formule des requêtes à un moteur de recherche, les requêtes pouvant contenir un ou plusieurs mots, voire un syntagme ou une phrase. Sachant que les taux de réussite de telles requêtes ne sont pas toujours satisfaisants, on se propose d'enrichir ces requêtes (en ajoutant des mots, ou en utilisant les opérateurs classiques de moteurs de recherche...) pour améliorer le résultat.</p> <p>Cet enrichissement sera réalisé sur la bases d'information de nature lexicale ou terminologique, stockées dans un réseau sémantique approprié.</p> <p>Le but du projet n'est pas de réaliser un tel réseau sémantique à grande couverture, mais de spécifier, pour un domaine conceptuel fixé, les heuristiques à utiliser et l'architecture du réseau correspondant. Le projet donnera lieu à la réalisation d'une maquette, au sens industriel du terme.</p>	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
	M1						
1999-2000	L3	Correcteur orthographique (approche graphémique)	<p>On se propose de réaliser un correcteur orthographique (lexical), qui, disposant d'un dictionnaire de formes fléchies, détecte les mots mal orthographiés, et propose si possible une correction.</p> <p>Le programme prend un texte (ASCII brut) en entrée, et pour chaque forme non présente dans le dictionnaire, propose à l'utilisateur :</p> <ul style="list-style-type: none"> • de choisir un de ces remplaçants <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences • de fournir un remplaçant <ul style="list-style-type: none"> o pour l'occurrence, ou pour toutes les occurrences o si le remplaçant n'est pas dans le lexique, ♣ de l'insérer dans le lexique pour ce texte, ♣ de l'insérer dans le lexique stable • d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique. <p>Les remplaçants proposés par le programme (premier cas plus haut) seront des mots figurant dans le dictionnaire, et ayant une certaine proximité avec le mot fautif.</p> <p>Pour la recherche des mots proches, on appliquera diverses heuristiques, basées sur la forme du mot :</p> <ul style="list-style-type: none"> • repérage des bigrammes ou trigrammes impossibles en français • redoublement ou dé-doublement de consonnes • suppression/insertion de diacritiques • insertion de lettres • interversion de lettres • etc. <p>Ces heuristiques, qui peuvent être plus nombreuses, seront étudiées linguistiquement afin de déterminer précisément leurs conditions d'application.</p>	Pascal ou C			
		Réaccentuation	<p>Le programme est chargé de remettre les accents et autres signes diacritiques manquant dans un texte fourni en typographie dite pauvre.</p> <p>On s'aidera d'un lexique de formes fléchies fourni. Bien sûr, certaines formes peuvent être réaccentuées directement, et d'autres sont ambiguës. Une fois que l'algorithme principal sera établi, on envisagera des heuristiques pour lever les ambiguïtés. Exemple :</p> <p>La ou le francais n'est pas accentue, il y a de la gene, mais quand le systeme m'accentue, je suis moins gene!</p> <p>On pourra s'inspirer du programme « Reacc » développé par le RALI à l'Université de Montréal.</p>	Pascal ou C			
		Correction d'étiquettes	<p>Disposant d'un corpus étiqueté et lemmatisé (chaque mot du corpus est accompagné de son étiquette grammaticale et du lemme lui correspondant [cf. figure 1]), il s'agit d'écrire un programme permettant de corriger d'éventuelles erreurs d'étiquetage.</p> <p>ce DETRMS ce paragraphe SUBSMS paragraphe résumé VER3B résumer les DETRMP les</p>	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>problèmes SUBMP problème qui PRELMP qui doivent VERB6 devoir être AUXE être résolus PPASMP résoudre pour PREP pour exploiter VINF exploiter des PDES des liaisons SUBSFP liaison de PDEA de télécommunicatio SUBSFS télécommunication par PREP par satellite SUBSMS satellite</p> <p>FIG 1. Exemple de phrase du corpus</p> <p>À ce programme, l'utilisateur fournira des règles qui auront la forme suivante :</p> <p>Changer l'étiquette Y1 du mot X en Y2 si les étiquettes Z1 et Z2 sont présentes dans le contexte de X, le contexte étant défini comme l'ensemble des cinq mots précédant et des cinq mots suivant X.</p> <p>Alors, le programme appliquera ces règles sur le corpus, en montrant à l'utilisateur les cas où la règle s'applique, de façon à ce qu'il puisse vérifier si la règle est pertinente.</p> <p>Exemple : au vu de l'étiquetage fautif</p> <p>I I DETRMS ' ' x usager usager SUBSMS désire désirer VERB3 établir établir VINF des des PDES liaisons liaison SUBSFP</p> <p>on pourra envisager une règle : L'étiquette PDES (=Y1) du mot des (=X) devient DETRMP (=Y2) si VERB3 (=Z1) et SUBSFP (=Z2) se trouvent dans le contexte du mot.</p> <p>On envisagera la possibilité d'appliquer plusieurs règles lors d'une passe, et on choisira de préférence de stocker ces règles dans un fichier. On notera que l'ordre d'application des règles est important.</p> <p>À titre d'exemple, vous proposerez un ensemble de règles qui permette de corriger les erreurs d'étiquetage sur ``des".</p> <p>On pourra envisager d'autres formes de règles.</p>				
		Calcul de collocations	<p>On veut réaliser un programme de calcul des collocations dans un corpus étiqueté et lemmatisé.</p> <p>Le but est de repérer et de compter les blocs de 2 ou 3 mots juxtaposés qui se retrouvent au moins deux fois dans le corpus. On s'intéressera aussi aux groupes « N à N », « N de N », etc., dont on calculera aussi les fréquences.</p> <p>La sortie du programme consiste donc en une liste, classée par fréquence, de toutes les collocations qui y apparaissent au moins deux fois.</p> <p>On attachera un soin particulier au choix de la structure de données, et à l'algorithme, qui doit faire le calcul en un temps raisonnable.</p>	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
		Extraction de sigles	<p>L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes) (par exemple EDF), et de leur forme développée (par exemple Electricité de France), utilisés dans un texte donné. Le programme aura deux fonctions :</p> <ol style="list-style-type: none"> 1. Création (extraction) de la liste (triée, normalisée) des sigles du texte donné. Attention aux formes diverses d'un même sigle (Inalf, I.N.A.L.F, INaLF). 2. Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte. <p>Le texte fourni est en ASCII brut.</p>	Pascal ou C			
		Justification et "hyphénation"	<p>Il s'agit de réaliser un programme qui justifie (au sens des traitements de texte) un texte fourni en ASCII, sur un nombre de colonnes donné, en découpant les mots selon les règles en usage pour le français.</p> <p>On supposera que tous les caractères ont la même dimension (« fonte proportionnelle »).</p> <p>Pour le découpage (éventuel) des mots (hyphenation), on prendra bien garde de distinguer les règles, qui seront stockées dans un fichier, du programme lui-même. On s'autorisera à utiliser un dictionnaire d'exceptions.</p>	Pascal ou C			
		Concordancier	<p>Réalisation d'un « concordancier » : disposant d'un corpus étiqueté et lemmatisé (chaque mot du corpus est accompagné de son étiquette grammaticale et du lemme lui correspondant (cf. figure 1), il s'agit d'extraire les phrases de ce corpus comportant des séquences d'étiquettes et/ou de lemmes données.</p> <p>Les séquences</p> <p>Les séquences sont des suites d'étiquettes et/ou de lemmes, et comportant éventuellement des symboles ayant une interprétation particulière (cf. expressions rationnelles) :</p> <p>*</p> <p>signifie zéro ou plusieurs occurrences de l'entité (étiquette grammaticale ou lemme) précédente,</p> <p>+</p> <p>signifie une ou plusieurs occurrences de l'entité précédente.</p> <p>Voici quelques exemples de séquences accompagnées de leur interprétation :</p> <p>SUBS ADJ un substantif, suivi d'un adjectif</p> <p>SUBS ADJ* PREP un substantif, suivi d'un nombre quelconque d'adjectifs, suivi d'une préposition</p> <p>SUBS ADJ+ PREP un substantif, suivi d'au moins un adjectif, suivi d'une préposition.</p> <p>Exemple de phrase à extraire</p> <p>Considérons la phrase de la figure 1. Cette phrase doit être extraite si la séquence recherchée est VINF des SUBS mais ne doit pas l'être si la séquence recherchée est SUBS ADJ des SUBS.</p> <p>Extraction des phrases</p> <p>Le programme propose à l'utilisateur de saisir une séquence (on rappelle que cette séquence peut comprendre des étiquettes et des lemmes), extrait les phrases qui comprennent cette séquence, et les stocke dans un fichier (on envisagera aussi leur visualisation). À titre d'exemple, si l'on s'intéresse au lemme des, on pourra demander au programme d'extraire du corpus les phrases correspondant aux séquences suivantes :</p> <ol style="list-style-type: none"> 1. VINF des SUBS 2. PREP VINF des ADJ* SUBS 3. PREP des ADV* ADJ* SUBS 4. SUBS ADJ+ des ADV* ADJ* SUBS <p>Remarques</p> <ul style="list-style-type: none"> • les étiquettes grammaticales utilisées dans les séquences sont plus générales que celles utilisées dans le corpus. En particulier, aucune information de genre et de nombre n'est donnée dans ces séquences, 	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<ul style="list-style-type: none"> • les prépositions ``à" et ``de" sont étiquetées PDEA dans le corpus sur lequel vous travaillerez mais doivent être considérées comme des prépositions. 				
	M1	Repérage des "entités nommées"	<p>On regroupe sous le terme « entités nommées » les noms de personnes, de lieux, les dates, etc. Le but du projet consiste à repérer automatiquement un sous-ensemble homogène parmi les entités (par exemple les noms de personnes ou les noms d'entreprises).</p> <p>Le programme envisagé procède par étapes successives sur le même texte. Dans une première étape, grâce aux dictionnaires spécialisés (cf. plus loin), on étiquette les mots simples. On partira d'un corpus déjà étiqueté morpho-syntaxiquement. La seconde étape, basée sur des règles, permettra de regrouper ces mots simples pour former des entités nommées. Enfin, une troisième étape, facultative dans ce projet, pourrait tenter de relier entre elles les entités co-référentes. Les éléments seront marqués dans le texte au moyen d'un système de balisage de type SGML.</p> <p>Le système repose donc, outre le programme lui-même, sur les ressources suivantes :</p> <ul style="list-style-type: none"> • des dictionnaires : <ul style="list-style-type: none"> o dictionnaire de la langue générale pour repérer les mots inconnus. o dictionnaires spécialisés de prénoms, de noms de sociétés connues, etc. o dictionnaires spécialisés d'« amorces », mots ou groupes de mots qui peuvent marquer le début d'une entité nommée. • une « grammaire » permettant de combiner des mots repérés de natures diverses (mots inconnus, prénoms, amorces...) pour en faire une entité (un nom de personne). <p>Les dictionnaires et la grammaire sont stockés dans des fichiers textes. On assure ainsi une séparation nette entre les données et les traitements (le programme).</p> <p>Exemple</p> <p>Soit le texte « M. Jack Lang a gagné la mairie. »</p> <p>"M." est stocké dans un dictionnaire "d'amorces"</p> <p>"Jack" est stocké dans un dictionnaire de prénoms</p> <p>"Lang" est inconnu du système</p> <p>La première étape permet d'obtenir :</p> <pre><TR_PERSON>M.</TR_PERSON> <PERSON>Jack</PERSON> <UFIRSTUNKNOWN>Lang</UFIRSTUNKNOWN> a gagné la mairie.</pre> <p>Une règle de la forme</p> <pre>TR_PERSON PERSON UFIRSTUNKNOWN ==> PERSON</pre> <p>va pouvoir être activée.</p> <p>D'où le résultat de la seconde étape :</p> <pre><PERSON>M. Jack Lang</PERSON> a gagné la mairie.</pre> <p>Une troisième étape permettrait d'identifier toutes les occurrences de Jack Lang</p> <pre><PERSON id=1>M. Jack Lang</PERSON> a gagné la mairie.</pre> <p>Sous réserve</p> <p>Je devrais pouvoir fournir une partie des dictionnaires à utiliser. Il est souhaitable que le système stocke dans un fichier les mots inconnus susceptibles d'être des entités, pour pouvoir améliorer rapidement et dynamiquement la couverture.</p>	Lisp, Prolog, Pascal, C, C++ ou java			
		Manipulation d'automates	On se propose de réaliser un programme, ou mieux, une bibliothèque de fonctions, qui permette de réaliser divers traitements sur des automates à nombre fini d'états.	Lisp, Prolog, Pascal, C, C++			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>On peut envisager une interface graphique (pour la sortie et même pour l'entrée des automates), mais on pourra aussi se contenter de travailler sur des tables de transition en Ascii, dans une notation à la PC-Kimmo. (Il faudra bien sûr envisager une structure de données interne pour travailler sur les automates.)</p> <p>Parmi les opérations que l'on devra pouvoir réaliser avec le système, il y aura au moins :</p> <ul style="list-style-type: none"> • Repérage de la catégorie de l'automate : complet ou non, déterministe ou non, avec ou sans epsilon-transitions... • Complétion d'un automate • Déterminisation d'un automate • Minimisation d'un automate • Suppression des epsilon-transitions d'un automate <p>Facultatif :</p> <ul style="list-style-type: none"> • Traduction d'une expression rationnelle en automate • Traduction d'un automate en expression rationnelle 	ou java			
		Réseaux sémantiques	<p>À partir de définitions provenant de dictionnaires électroniques, préalablement étiquetées, il s'agit de construire un « réseau sémantique ».</p> <p>Pour chaque entrée (qui peut correspondre à un des sens d'une lexie), on repèrera les mots pleins qui participent à sa définition (il faut donc (1) distinguer les mots pleins des mots « outils », (2) repérer la définition proprement dite parmi l'ensemble des informations associées à une entrée (catégorie, exemples, synonymes...)), et on construira un réseau reliant l'entrée à tous ces mots pleins.</p> <p>Ce réseau pourra être stocké sous forme de fichier Ascii, dans un format du genre de celui de WordNet.</p> <p>On pourra aussi, éventuellement, proposer une interface graphique permettant de visualiser graphiquement le réseau, voire de le modifier.</p> <p>Dans un deuxième temps, on réalisera un programme exploitant ce réseau pour désambiguïser les sens d'un mot en contexte. Principe : étant donnée une phrase contenant le mot concerné, on repère les mots pleins qui l'entourent (contexte), et on recherche ces mots dans le réseau. Alors le sens correspondant est celui qui est le plus proche (dans un sens qu'il faut précisément définir) des mots pleins de son contexte.</p>	Lisp, Prolog, Pascal, C, C++ ou java			
		Génération de texte en LFG (lisp)	<p>Le programme doit pouvoir afficher la ou les phrases correspondant à une structure fonctionnelle supposée bien formée (cohérente et complète). Le programmeur fournira une petite grammaire LFG pour la langue de son choix.</p> <p>1. Première étape : Écrire un programme qui génère le langage correspondant à une grammaire CF quelconque. Le programme pourra « boguer » en cas de récursion gauche. On pourra produire par exemple un réseau d'automates.</p> <p>2. Seconde étape : Compléter ce programme pour tenir compte des équations fonctionnelles. Seul le langage correspondant à la structure fonctionnelle devra être produit. Procéder par étapes :</p> <ul style="list-style-type: none"> o analyser les équations du type $v = (\wedge \text{obj})$ o ensuite les equations du type $v = (\wedge (v \text{ pcas }) - \text{obj})$ o ensuite les équation contraintes o ensuite les equations existentielles <p>3. Troisième étape : Retrouver les formes fléchies dans un dictionnaire à partir des lemmes et informations morphologiques des mots.</p>	Lisp, Prolog, Pascal, C, C++ ou java			
		"Complétion" automatique	<p>On se propose de réaliser un programme qui accélère la vitesse de frappe d'un utilisateur en proposant - de façon intelligente - la fin des mots en cours de frappe.</p> <p>Un tel système, pour être vraiment performant, devrait faire de la prédiction syntaxique en temps réel. On se contentera dans ce projet de faire un système qui propose la suite la plus probable des mots.</p>	Emacs Lisp			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>La méthode repose sur un calcul de probabilité basé sur la fréquence d'apparition de chaque mot dans un corpus préalablement fourni. On construit un graphe dont chaque chemin correspond à un mot (une lettre par arc), et dont chaque arc porte une probabilité.</p> <p>Par exemple, la probabilité de la prochaine lettre de la suite ``v-o-l-u'' est respectivement de 0/9, 8/9 (3+5/9), 1/9, 0/9 pour les lettres `b', `m', `p' et `t' si nous avons dans le corpus :</p> <p>0 occurrences du mot volubile 3 " volume 5 " volumineux 1 " voluptueux 0 " volute</p> <p>La lettre la plus probable, ici ``m'', sera proposée à l'utilisateur, et un nouveau calcul est effectué pour la suite de lettres ``v-o-l-u-m''.</p> <p>Un soin particulier sera apporté à l'interface de ce programme : l'utilisateur doit pouvoir accepter ou revenir facilement sur les propositions du système, etc. Un langage particulièrement adapté au projet serait Emacs Lisp, ce qui permettrait d'intégrer la complétion à l'éditeur Emacs.</p>				
		Désambiguïsation par automates	<p>Partant d'un texte dont chaque mot est étiqueté avec toutes ses catégories morpho-syntaxiques possibles, on va s'attacher à lever certaines ambiguïtés non pertinentes, c'est-à-dire à supprimer certaines des étiquettes attachées à certains mots.</p> <p>Pour mener à bien cette tâche, on extraiera du corpus donné les éléments pertinents, ce qui permettra d'établir des « règles contextuelles », que l'on appliquera ensuite au corpus.</p> <p>1. Étude linguistique</p> <p>On s'intéressera à une catégorie donnée, par exemple « verbe conjugué ». On extraira du texte une première liste de mots qui peuvent être des verbes conjugués (ceux dont la liste d'étiquettes comprend la catégorie V), puis, après désambiguïsation, une deuxième liste comprenant seulement les mots qui sont effectivement des verbes conjugués dans le texte.</p> <p>On relèvera alors les contextes locaux (catégories précédentes, catégories suivantes) qui permettent de choisir la catégorie retenue pour un mot donné en français.</p> <p>Par exemple, si X est marqué N ou V, et si X est précédé d'un Det, alors X doit être marqué N (l'étiquette V doit être ôtée).</p> <p>Etablir ainsi des séquences de 2 ou 3 catégories impossibles en français.</p> <p>On aura intérêt à subdiviser certaines catégories du DELAF, par exemple mettre à part parmi les pronoms, les formes clitiques (je, tu, il(s), elle(s), nous, vous, se, on, ce, y, en, le, la, les, lui, leur).</p> <p>Attention au cas où les mots suivants ou précédents ont plusieurs catégories possibles.</p> <p>2. Programmation et implémentation</p> <p>Une fois les règles dégagées, on réalisera un programme qui procédera à la désambiguïsation du corpus en appliquant les règles. On conseille d'implémenter les règles sous forme de transducteur.</p> <p>On fournit deux textes étiquetés (d'une dizaine de pages) à l'aide du DELAF, qui serviront de base à l'élaboration des règles de désambiguïsation.</p> <p>Notations du DELAF :</p> <p>Général Pour les verbes V verbe N nom DET déterminant (article) DETP déterminant possessif DETQ déterminant numéral</p>	Lisp, Prolog, Pascal, C, C++ ou java			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			INTJ interjection CNJS conjonction de subordination CNJC conjonction de coordination A adjectif PRON pronom ADV adverbe Xin préfixe P présent indicatif S subjonctif présent Y impératif K participe passé F futur W infinitif I imparfait C conditionnel G participe présent T subjonctif imparfait Inv invariable 1, 2 ou 3 : personne m, f, s, p : genre, nombre Exemple Forme générale : mot lemme cat:morpho ... vient venir V33:P3s livre livre N1:ms livre N21:fs livrer V3:P1s:P3s:S1s:S3s:Y2s				
		Compression de texte	On doit réaliser un couple de programmes qui permette de compresser-décompresser un texte. Le principe de l'algorithme de compression est le suivant : sur des bases statistiques, en fonction du texte en entrée, on détermine des chaînes les plus longues possible (en particulier préfixes et suffixes) présentes avec une fréquence raisonnable, et on remplace dans le texte toutes les occurrences de ces chaînes par des non-mots -ou des mots absents du texte- (p. ex. « aa »), plus courts. Bien sûr, le fichier compressé contiendra aussi une table de conversion, qui permettra au programme de décompression de reconstituer le fichier original.	Lisp, Prolog, Pascal, C, C++ ou java			
1998-1999	L3						
	M1						
1997-1998	L3	Extraction de phrases		Pascal ou C			
		Correcteur orthographe	On veut réaliser un correcteur orthographique (lexical) qui fait des propositions de correction lorsqu'il rencontre un mot n'appartenant à son lexique. Fonctionnement : on dispose d'un lexique de formes fléchies. Un texte est fourni en entrée du programme (ascii brut, sans marquage). Pour chaque mot du texte non présent dans le lexique, l'algorithme déterminera un ou plusieurs mots proches (mais figurant dans le lexique), et proposera à l'utilisateur :	Pascal ou C			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>de choisir un de ces remplaçants pour l'occurrence, ou pour toutes les occurrences</p> <p>de fournir un remplaçant pour l'occurrence, ou pour toutes les occurrences</p> <p>si le remplaçant n'est pas dans le lexique, de l'insérer dans le lexique pour ce texte, de l'insérer dans le lexique stable</p> <p>d'ignorer la correction, c'est-à-dire conserver le mot initial, pour cette occurrence ou pour toutes, en l'insérant ou non dans le lexique.</p> <p>En qui concerne la détermination de la proximité entre deux mots, on réalisera un « panachage » entre les deux méthodes classiques :</p> <p>approche graphémique : règles basées sur la forme : redoublement de consonnes, substitutions de voyelles, interversion de lettres (voisines sur le clavier, p.ex), etc.</p> <p>approche phonétique : règles basée sur la proximité phonétique : substitution de `k' par `qu', `c' (si suivi de la bonne voyelle), de `ph' par `f', gestion des `e' muets, etc.</p> <p>Remarques</p> <p>Il n'est pas raisonnable de charger tout le dictionnaire en mémoire. Il faut envisager des chargements partiels.</p> <p>On peut envisager un dictionnaire parallèle de formes phonétiques associées, ou une table de correspondance graphème-phonème. Dans les deux cas, ce n'est pas incompatible avec le développement d'heuristiques spécifiques.</p> <p>On peut limiter la complexité des calculs en tenant compte des bigrammes possibles et impossibles en français.</p>				
		Correction d'étiquettes		Pascal ou C			
		Prédiction d'étiquettes	<p>Disposant d'un lexique de formes fléchies, associées à une étiquette (partie du discours, + informations morpho-syntaxiques), et d'un texte à étiqueter, le programme doit étiqueter les formes reconnues, et tenter de deviner, uniquement sur des bases morphologiques, l'étiquette probable des formes qui ne sont pas dans le lexique.</p> <p>On se servira d'une table classique morphologique (grammaire Becherelle), la terminaison du mot appartient à l'une des terminaisons décrites, on en prédit l'étiquette.</p> <p>Exemple : la terminaison de vermifugerions "erions" donne à croire que le mot est un verbe au conditionnel présent première personne du singulier.</p>	Pascal ou C			
		Extraction de sigles	<p>L'objectif est d'être capable de construire, semi-automatiquement, la liste des sigles (acronymes), et de leur forme développée, utilisés dans un texte donné. Le programme aura trois fonctions :</p> <p>Création (extraction) de la liste (triée, normalisée) des sigles du texte donné.</p> <p>Correction (automatique) si les sigles apparaissent sous plusieurs formes dans le texte (I.N.A.L.F, Inalf, INaLF).</p> <p>Proposition (interactive) de syntagmes candidats pour la forme développée, quand on en trouve dans le texte.</p> <p>Le texte fourni est en ASCII brut.</p>	Pascal ou C			
	M1	Complétion automatique		Lisp, Prolog, Pascal, C, C++ ou java			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
		Désambiguïsation par automates		Lisp, Prolog, Pascal, C, C++ ou java			
		Traduction de nombres		Lisp, Prolog, Pascal, C, C++ ou java			
		Compréhension de dates		Lisp, Prolog, Pascal, C, C++ ou java			
		Analyseur syntaxique		Lisp, Prolog, Pascal, C, C++ ou java			
1996-1997	L3		<p>Écriture d'un programme permettant :</p> <ol style="list-style-type: none"> d'extraire des phrases correspondant à des séquences données, de corriger des étiquettes grammaticales à partir de règles à définir. <p>1. Extraction de phrases correspondant à des séquences données</p> <p>Disposant d'un corpus étiqueté et lemmatisé (chaque mot du corpus est accompagné de son étiquette grammaticale et du lemme lui correspondant), il s'agit, dans un premier temps, d'extraire les phrases de ce corpus comportant des séquences d'étiquettes et/ou de lemmes données.</p> <p>1.1. Les séquences</p> <p>Les séquences représentent des suites admissibles d'étiquettes grammaticales et/ou de lemmes que l'on peut représenter sous forme d'expressions régulières.</p> <p>Deux symboles ont une signification particulière dans les séquences:</p> <p>*</p> <p>signifie zéro ou plusieurs occurrences de l'entité (étiquette grammaticale ou lemme) précédente,</p> <p>+</p> <p>signifie une ou plusieurs occurrences de l'entité précédente.</p> <p>Voici quelques exemples de séquences accompagnées de leur interprétation :</p> <p>SUBS ADJ un substantif, suivi d'un adjectif</p> <p>SUBS ADJ* PREP un substantif, suivi d'un nombre quelconque d'adjectifs, suivi d'une préposition</p> <p>SUBS ADJ+ PREP un substantif, suivi d'au moins un adjectif, suivi d'une préposition.</p> <p>1.2. Exemple de phrase à extraire</p> <p>Considérons la phrase suivante extraite du corpus :</p> <p>ce DETRMS ce</p> <p>paragraphe SUBSMS paragraphe</p> <p>résumé VER3B résumer</p> <p>les DETRMP les</p> <p>problèmes SUBMP problème</p> <p>qui PRELMP qui</p> <p>doivent VERB6 devoir</p> <p>être AUXE être</p>	Pascal			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>résolus PPASMP résoudre pour PREP pour exploiter VINF exploiter des PDES des liaisons SUBSFP liaison de PDEA de télécommunicatio SUBSFS télécommunication par PREP par satellite SUBSMS satellite</p> <p>Cette phrase doit être extraite si la séquence recherchée est VINF des SUBS mais ne doit pas l'être si la séquence recherchée est SUBS ADJ des SUBS.</p> <p>1.3. Extraction des phrases</p> <p>Il vous est ici demandé d'extraire du corpus les phrases correspondant aux séquences suivantes :</p> <ol style="list-style-type: none"> 1. VINF des SUBS 2. PREP VINF des ADJ* SUBS 3. PREP des ADV* ADJ* SUBS 4. SUBS ADJ+ des ADV* ADJ* SUBS <p>A chaque séquence sera associé un fichier dans lequel vous rangerez les phrases correspondant à la séquence.</p> <p>Remarques</p> <ul style="list-style-type: none"> • les étiquettes grammaticales utilisées dans les séquences sont plus générales que celles utilisées dans le corpus. En particulier, aucune information de genre et de nombre n'est donnée dans ces séquences, • les prépositions ``à" et ``de" sont étiquetées PDEA dans le corpus sur lequel vous travaillerez mais doivent être considérées comme des prépositions. <p>Indication Il est possible de représenter les séquences sous forme de tableau où chaque élément correspond à une entité. La séquence SUBS ADJ+ PREP peut ainsi être représentée par le tableau suivant, dans lequel l'élément ``+" indique qu'on peut répéter l'élément précédent. Une séquence est reconnue dans une phrase lorsqu'on a réussi à parcourir tout le tableau en lisant la phrase.</p> <p>SUBS ADJ + PREP</p> <p>2. Correction des étiquettes grammaticales</p> <p>Il s'agit ici d'écrire des règles et un programme permettant de corriger les éventuelles erreurs d'étiquetage que vous aurez pu relever sur le mot ``des".</p> <p>Les règles auront la forme suivante :</p> <p>Changer l'étiquette Y1 du mot X en Y2 si les étiquettes Z1 et Z2 sont présentes dans le contexte de X, le contexte étant défini comme l'ensemble des cinq mots précédant et des cinq mots suivant X.</p> <p>Vous pourrez ajouter des contraintes sur ces règles (comme par exemple le fait que Z1 et Z2 doivent se suivre).</p> <p>Il vous faudra vérifier que les règles que vous écrivez ne génèrent pas plus d'erreurs qu'elles n'en corrigent.</p> <p>3. Extension</p> <p>Vous pourrez définir de nouvelles séquences associées à ``des" si vous le jugez utile (ou si vous avez terminé le projet en avance). Vous pourrez également vous intéresser à d'autres mots que ``des".</p>				
			Écriture d'un programme permettant :	Pascal			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			<p>1. d'extraire des pseudo-affixes en vue de sélectionner les séquences à abrégger, 2. de générer automatiquement des codes correspondant aux séquences retenues, 3. de prendre en compte ces codes pour simplifier l'écriture d'un mot.</p> <p>1. Extraction de pseudo-affixes</p> <p>Il s'agit ici de trouver, à partir des mots présents dans un corpus, les séquences de caractères les plus représentées. Seules les séquences supérieures à n (à choisir et/ou déterminer) caractères seront considérées. Il faudra classer les séquences obtenues par ordre décroissant du nombre d'occurrences.</p> <p>L'extraction reposera sur l'algorithme présenté dans l'article de Nagao et Mori ci-joint. Vous commencerez par vous intéresser aux séquences finales, qu'on pourrait qualifier de pseudo-suffixes. Le corpus sur lequel vous travaillerez est un corpus étiqueté et lemmatisé identique à celui décrit dans le projet 1.</p> <p>Une fois les séquences extraites, vous devrez en choisir un certain nombre sur lesquelles vous travaillerez, ce choix devant être motivé.</p> <p>Remarque Il vous faudra choisir entre un travail sur formes fléchies et sur lemmes pour extraire vos séquences. Vous pouvez choisir l'un ou l'autre (les deux présentent un intérêt), mais vous devez indiquer les raisons de votre choix.</p> <p>2. Génération automatique de codes</p> <p>A partir des séquences retenues, il vous est ici demandé de générer automatiquement un code pour chacune d'elles. Ce code prendra la forme d'une suite de ``a" et de ``b". La séquence ``ment" peut par exemple être représentée par la suite ``aa", de telle sorte que pour signifier ``changement" on peut se contenter de taper ``changeaa". Le gain en lettres est ici faible, mais il faut tenir compte du fait que le code généré doit être non ambigu. Ainsi, si l'on avait choisi ``a" pour représenter la séquence ``ment", on n'aurait pu décider si ``changea" représentait ``changement" ou le passé simple du verbe ``changer".</p> <p>Vous pourrez vous inspirer de l'algorithme de Huffman pour générer les codes. Je vous fournirai cet algorithme en temps utile.</p> <p>3. Prise en compte des codes</p> <p>Cette partie concerne le remplacement automatique des codes par les séquences associées. Il s'agit ici d'écrire une fonction lisant un mot, fourni par l'utilisateur, éventuellement incomplet mais comportant un code et renvoyant le mot complet. Suivant le code généré précédemment, plusieurs stratégies peuvent être utilisées. Si le code reste ambigu, on peut demander à l'utilisateur de séparer le mot incomplet du code par un caractère de contrôle (``?", par exemple, qui n'est attesté à l'intérieur d'aucun mot de la langue française). Si le code généré est non ambigu, un tel caractère est superflu.</p> <p>4. Extension</p> <p>On peut envisager le même traitement sur les séquences initiales des mots (pseudo-préfixes), voire sur des séquences quelconques. La non ambiguïté du code devient alors un critère de première importance. Essayer d'envisager un tel traitement.</p>				
	M1	Traduction de nombres écrits en lettres	<p>1. Etude linguistique</p> <p>Etablir vocabulaire et grammaire (catégories et graphe de transition) pour les nombres dans chacune des langues choisies.</p> <p>2. Programme et implémentation</p> <p>Implémenter un automate (par exemple sous forme de transducteur) qui reconnaisse si la séquence (entrée par l'utilisateur) est bien un nombre dans la langue source, et la traduise en chiffre.</p> <p>Exemples : mille deux cents -> 1200</p> <p>un mille -> erreur</p> <p>mille un -> 1001</p> <p>Implémenter un automate qui génère toutes les séquences possibles dans la langue cible.</p>	Emacs Lisp ou Prolog			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			Exemples : 1200 -> one thousand (and) two hundred, a thousand (and) two hundred mille un -> one thousand and one, a thousand and one				
		Compréhension de dates écrites en lettres	<p>1. Etude linguistique</p> <p>Etablir vocabulaire et grammaire (catégories et graphe de transition) pour les dates dans la langue choisie.</p> <p>2. Programmation et implémentation</p> <p>Soit analyser la séquence linguistique (entrée au clavier par l'utilisateur), vérifier que c'est bien une date et la traduire en chiffre. Attention aux expressions elliptiques, et aux déictiques.</p> <p>Exemples : le dix-huit février -> 18/02/1997</p> <p>le décembre 18 -> ceci n'est pas une date</p> <p>demain -> 23/02/1995 (appel horloge système)</p> <p>le lundi 14 décembre 36 -> 14/12/1936</p> <p>en décembre 18 -> 01/12/1918 < d < 31/12/1918</p> <p>Soit générer toutes les séquences linguistiques a partir d'une date en chiffres (entrée au clavier par l'utilisateur) :</p> <p>Exemples : 18/02/95 -> le dix-huit février mille-neuf-cent-quatre-vingt-quinze, le 18 février 1995, etc. (Calculs supplémentaires pour : le samedi 18 février, le troisième samedi du mois courant...)</p> <p>23/02/95 -> le jeudi 23 février, demain (appel horloge système)...</p> <p>29/02/95 -> 1995 n'est pas une année bissextile</p> <p>32/06/97 -> Ceci n'est pas une date</p>	Emacs Lisp ou Prolog			
		Désambiguation par automates et extraction de catégories à partir d'un texte étiqueté	<p>Il s'agit, à partir d'un texte découpé en mots, chacun étiqueté avec toutes ses catégories grammaticales possibles, de lever certaines ambiguïtés non pertinentes.</p> <p>On s'attachera à l'extraction d'une catégorie donnée, par exemple ((Verbe conjugué)). On extraira du texte une première liste de mots qui peuvent être des verbes conjugués, puis, après désambiguation, une deuxième liste comprenant seulement les mots qui sont effectivement des verbes conjugués dans le texte.</p> <p>Il est conseillé de procéder comme suit :</p> <p>1. Etude linguistique</p> <p>Relever les contextes locaux (catégories précédentes, catégories suivantes) qui permettent de choisir la catégorie retenue pour un mot donné en français.</p> <p>Par exemple, si X est marqué N ou V, Det + X => X = N.</p> <p>Etablir ainsi des séquences de 2 ou 3 catégories impossibles en français</p> <p>On aura intérêt à subdiviser certaines catégories du DELAF, par exemple mettre à part parmi les pronoms, les formes clitiques (je, tu, il(s), elle(s), nous, vous, se, on, ce, y, en, le, la, les, lui, leur).</p> <p>Attention au cas où les mots suivants ou précédents ont plusieurs catégories possibles.</p> <p>2. Programmation et implémentation</p> <p>On conseille d'implémenter les règles dégagées sous forme de transducteurs.</p> <p>On fournit deux textes étiquetés (d'une dizaine de pages) à l'aide du DELAF, qui serviront de base à l'élaboration des règles de désambiguïsation. Le programme devra être testé sur ces deux textes, ainsi que sur un troisième qui n'aura pas été étudié a priori.</p> <p>Rappel des notations du DELAF :</p> <p>V : verbe</p> <p>N : nom</p>	C ou Pascal			

ANNEE	NIVEAU	SUJET	DESCRIPTIF	LANGAGE	RESPONSABLE	DIFFICULTE	GROUPE
			DET : déterminant (article) DETP : déterminant possessif DETQ : déterminant numéral INTJ : interjection CNJS : conjonction de subordination CNJC : conjonction de coordination A : adjectif PRON : pronom ADV : adverbe Xin : préfixe Pour les verbes : P : présent indicatif S : subjonctif présent Y : impératif K : participe passé F : futur W : infinitif I : imparfait C : conditionnel G : participe présent T : subjonctif imparfait Inv : invariable 1, 2 ou 3 : personne m, f, s, p : genre, nombre Exemple Forme générale: mot lemme cat:morpho ... vient venir V33:P3s livre livre N1:ms livre N21:fs livrer V3:P1s:P3s:S1s:S3s:Y2s				
				Emacs Lisp ou Prolog			